

The AI to Read Your Mind: Vetting Response Integrity from cross-Talk dependencies in Adversarial Surveys

Nicholas Sizemore¹, Royce Lee², James Evans^{3,4}, Robert Gibbons^{1,4,5,6}, and Ishanu Chattopadhyay^{1,4,6}★

¹Department of Medicine, University of Chicago, Chicago, IL 60637, USA

²Department of Psychiatry, University of Chicago, Chicago, IL 60637, USA

³Department of Sociology, University of Chicago, Chicago, IL 60637, USA

⁴Committee on Quantitative Methods in Social, Behavioral, and Health Sciences, University of Chicago, Chicago, IL 60637, USA

⁵Department of Public Health Sciences, University of Chicago, Chicago, IL 60637, USA

⁶Center for Health Statistics, University of Chicago, Chicago, IL 60637, USA

★To whom correspondence should be addressed: e-mail: ishanu@uchicago.edu.

Abstract: Exaggerating or faking the symptoms of a mental health disorder can confound structured diagnostic interviews and hinder clinical psychiatric assessments^{1,2}. We introduce an artificial intelligence (AI) framework for detecting symptom fabrication in mental health assessments, illustrated here for Post-Traumatic Stress Disorder (PTSD) diagnoses, for which malingering is a known problem^{3,4}, partially ascribable to the potential for secondary financial gain from positive diagnoses. Algorithm **veRITAS** employs novel generative AI to infer statistical dependencies inherent in true response patterns, and flags responses which violate these subtle constraints. With a study sample of $n = 651$ patients, **veRITAS** is estimated to have a Area Under the Curve (AUC) of $\geq 0.95 \pm 0.02$, with sensitivity $> 95\%$, specificity $> 88\%$ respectively, and positive likelihood ratios between $9.9 - 19.77$ achievable based on the population prevalence of malingering in the context of PTSD diagnosis. We show that in our methodology having training in forensic psychiatry, or other relevant mental health experience, is not helpful in deceiving the algorithm. Our tool offers an objective, disease-specific, fast (average time ≤ 4 min) approach to detect fake PTSD, and if adopted, can ensure that healthcare resources and disability concessions reach those genuinely in need, while helping to maintain integrity of clinical data. Moreover, the ability to identify and help patients who might be malingering due to other mental health conditions, poverty or socio-economic compulsions can improve general health outcomes in disadvantaged communities.

ONE SENTENCE SUMMARY

Generative AI-based lie-detector flags fake mental health symptoms with 95% accuracy immune to deception even by trained psychiatrists.

INTRODUCTION

Diagnosis and measurement of severity of mental health disorders typically rely on structured interviews⁵, clinician symptom ratings, or patient self-reports, and are potentially susceptible to intentional fabrication of symptoms^{1,2}, referred to as “malingering”. In the context of the measurement and diagnosis of Post-Traumatic Stress Disorder (PTSD) as an example, subjective stressors coupled similarity in presentation, relatively easy access to information on how to fake PTSD, and financial incentives related to a positive diagnosis in legal procedures or disability claims, is known to hinder accurate clinical assessments. While there are existing elementary approaches to construct “lie scales” dating back to the Minnesota Multiphasic Personality Inventor (MMPI)⁶, here, we present an approach based on a novel generative artificial intelligence (AI) model to flag malingering suspects rapidly, efficiently and accurately.

We illustrate use of our methodology using symptom-level data for PTSD. Clinically, PTSD is an anxiety disorder that can develop after experiencing a traumatic event. In the United States, disability compensation is available for those with mental health disorders, which while being a crucial resource for the truly afflicted, incentivizes

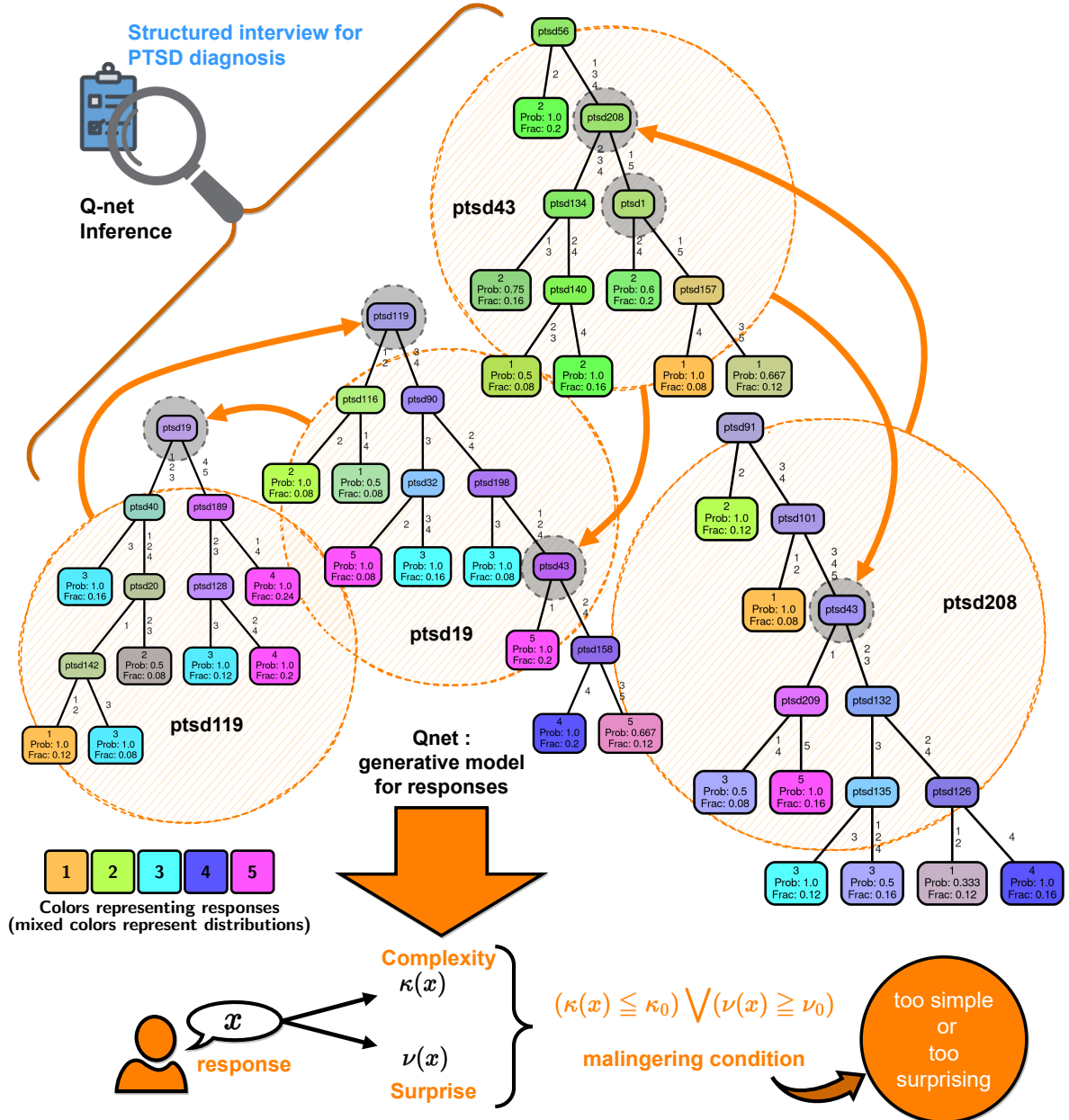


Fig. 1: Conceptual framework. Using a dataset of responses to a validated structured interview for PTSD diagnoses, along with physician-validated clinical diagnoses, we infer generative models for responses for PTSD patients. In our framework for detecting malingering, we flag responses as those which are highly “surprising” (defined as violating inferred cross-dependencies between individual response items) or are too simple (lacking complexity in the response patterns typical of non-malingered responses). The precise “malingering condition” shown above is validated from theoretical considerations as well as field experimental data.

malingering^{3,4}. Thus, faking PTSD to access medical treatment, commit insurance, personal injury and other frauds, or in an attempt to evade criminal liability and penalties^{7–10} is unfortunately not rare. While PTSD is a serious mental health condition associated with substance use disorder, mood disorder, anxiety disorder, personality disorder, increased morbidity, and possibly with increased mortality^{11,12}, false diagnoses and fabrication of self-reported symptom severity can cause substantial financial drain^{13,14} to healthcare systems, divert crucial resources from where they are needed¹⁵, and interfere with study outcomes by introducing inaccuracies in clinical data¹⁶. Accurate disambiguation of true and fabricated PTSD is therefore of high importance, especially with sources suggesting that over 20% of personal injury cases, as well as 20% of the Veterans seeking combat compensation could be fabricating their condition^{3,17–19}. In addition to curbing financial fraud, ability to flag such incongruities can help address the myriad of factors that can drive malingering behavior, including other mental health conditions, a lack of access to healthcare^{20,21} and an inability to seek help arising from poverty and other broad-ranging socio-economic conditions.

Despite the general difficulty in formulating principles to detect malingering²², multiple standardized tests^{23,24} and

TABLE 1: Demographic Characteristics and Malingering Success Rates (PL Dataset)* at 94.2% sensitivity

characteristics	mean Completion Time [s]	mean age [years]	no. of participants	malingering success rate (%)
Race: Asian	188.4	34.3	24	8.3
Race: Black	329.5	38.0	33	18.2
Race: Mixed	195.8	31.4	20	5.0
Race: Other	286.2	38.5	13	0.0
Race: White	186.5	42.7	220	4.1
Sex: Female	201.5	40.5	167	4.8
Sex: Male	213.5	41.0	141	7.1
Residence: United Kingdom	213.5	43.1	110	4.5
Residence: United States	202.9	39.4	200	6.5
All participants	206.6	40.7	310	5.8

*Using $\kappa_0 = 1$, $\nu_0 = 0.76$, $\mu_0 = 1.35$.

TABLE 2: Lower Bounds on Performance trade-offs at Different Population Prevalences of PTSD Malingering†

prev.	sensitivity	specificity	ppv	npv	acc	LR+	LR-
0.15	0.930±0.014	0.95	0.777±0.003	0.987±0.002	0.949±0.002	19.77±0.347	0.074±0.015
0.20	0.942±0.012	0.93	0.831±0.002	0.984±0.003	0.948±0.003	19.77±0.347	0.074±0.015
0.25	0.948±0.012	0.91	0.868±0.002	0.981±0.004	0.947±0.003	19.77±0.347	0.074±0.015
0.30	0.956±0.010	0.88	0.894±0.001	0.979±0.004	0.945±0.004	9.988±0.890	0.074±0.015

† Abbreviations: Population prevalence (prev.), Positive Predictive Value (ppv), Negative Predictive Value (npv), Accuracy (acc), Positive Likelihood Ratio (LR+), Negative Likelihood Ratio (LR-). 99% confidence intervals calculated for $n = 310 + 0.6 \times 304 = 492$.

validity assessment tools²⁵ have been proposed, with limited success. These tools typically aim to incorporate patterns observed in diagnostic populations that might disambiguate fabricated symptoms from real ones or ask similar or related questions multiple times to verify consistency. However, existing approaches do not target specific disorders, almost always require expert interpretation, are subjective, and by design are unlikely to be effective against a malingerer with psychiatric training (See Table 4). Other strategies with physiological monitoring and linguistic analysis^{26–29} cannot be easily adopted in structured or semi-structured interviews or patient reported outcome (PRO) measurement.

In *VeRITAS* we leverage the fact that both clinician-rated and patient self-reported symptoms have statistical dependencies arising from the nature of the questions themselves (Fig. 1) and are modulated by the latent trait or condition we are attempting to measure *e.g.*, PTSD diagnosis and/or severity. We operationalize this principle without requiring human-understanding of the specific items being administered; thus, making the approach specific to the disease at hand (PTSD), while being potentially generalizable to other disorders if appropriate training data are available.

Our key finding is that the subtle cross-dependencies between the symptom items are challenging to mimic on-the-fly, even with training in forensic psychiatry. In particular, maintaining the right amount of expected structure in a sequence of responses to a structured diagnostic interview, measurable via our “complexity” and “surprise” parameters proves is difficult. Thus, *VeRITAS* offers a robust approach for flagging malingering subjects, can be trained to target specific disorders, can be administered in less than 4 minutes on average, and requires no subjective interpretation.

RESULTS

Participants and Data Sources

Our first dataset (referred to as the VA dataset) comprises $n = 304$ participants recruited at a Veterans Health Administration facility for an earlier study³⁰. Veterans between the ages of 18 and 89 years were recruited with written informed consent. Once eligibility was determined by the study team, participants completed a PTSD-symptom questionnaire from the CAT-PTSD item bank^{30,31}, comprising 211 items, including some items from the PTSD Checklist (PCL-5). Participants were also interviewed using the Clinician-Administered Scale for PTSD for Diagnostic and Statistical Manual of Mental Disorders, Fifth Edition (DSM-5), which resulted in 86 PTSD diagnoses, and 218 participants deemed as not having PTSD. We used 60% of the data for inferring our generative models, and the rest were used for validation, including determining the out-of-sample AUC for identifying PTSD vs no-PTSD cases. The possibility of malingering was not recorded in this dataset.

Our second dataset (the PL dataset) comprises results from online surveys conducted by a third party vendor

TABLE 3: Demographic characteristics of participating mental health professionals

Area of Expertise	
Neuropsychology	11
Forensic Psychiatry	8
Psychology	4
Forensic Neuropsychology	3
Neurology	1
Sex	
Male	14
Female	13
Primary Institution of Affiliation	
Northwestern University	21
University of Illinois Chicago	3
NorthShore University HealthSystem	1
University of Chicago Medicine	1
Rush University Medical Center	1

(Prolific) hired by the study team. Referred to here as the Prolific dataset, it comprises $n = 310$ participants (200 in the US, and 110 in the UK), screened for the absence of past or present mental health diagnoses. The participants were asked to fake symptoms of PTSD while taking the `VeRITAS` interview.

Our third dataset (the PS cohort) comprises $n = 27$ mental health professionals including forensic psychiatrists recruited from academic institutions in the Chicago region (See Table 3) who, akin to the PL participants, took the `VeRITAS` test aiming to malingering and get a false PTSD diagnosis. The objective in collecting the PS cohort was to test if extensive knowledge of the relevant mental health pathologies makes it easier to defeat the test.

In total, we considered $n = 651$ participants comprising US Veterans, general population in the US and the UK, and mental health professionals, with broad representation across sexes, race and ethnicity. Average completion times and detailed demographic composition of the respondents in the PL dataset are shown in Table 1, and that for the PS cohort are shown in Table 3. Both the PL and the PS datasets were generated via a “web-app” implementation of the `VeRITAS`. Unlike the PL dataset for which identifying respondent information is available to the third-party vendor (but not released to the study team), for the PS cohort we only collected de-identified responses, and hence, while we knew who were in the complete set of respondents for the PS cohort, we did not collect information to identify which response was generated by which individual.

Principle of Characterizing Malingered Responses

Our first insight behind `VeRITAS` is that response sequences from malingering subjects necessarily have high average “surprise”, *i.e.* deviate more on average from the context-specific model predictions of symptom ratings (*i.e.* item responses). Here context-specificity refers to the dependence of a response to responses to other items, which might be indicative of behavioral or mental health phenotypes. Our other insight, based on observations, is that attempts to mimic true PTSD tends to generate “over-structure” compared to responses from participants with actual PTSD, *i.e.*, malingering tends to manifest too much regularity in the response patterns. In other words, mimicked responses are less “complex”. Here, we understand complexity in the formal sense of Kolmogorov³²: more complex objects are less compressible, and perfectly random sequences being not significantly compressible at all. Thus, a sequence of all 1’s (*i.e.* symptom absence) is very compressible (since, instead of storing individual responses, one could just remember such responses as “all 1’s”), but a sequence generated from the sequential tosses of a fair coin is not very compressible, *i.e.* has less structure. We hypothesize that true response sequences tends to have maximal randomness, conditional on being constrained by the emergent cross-dependencies. In other words, true responses should have just the right amount of identifiable structure, and *no more*. More structure will make it less random, and less or inappropriate structure will increase surprise. Thus too much, too little or deviant structures are all indicative of deviant responses, interpreted here as probably malingering. It is important to note that these criteria (upper bound on surprise and lower bound in complexity) are based solely on algorithmic properties of the response sequence, and do not require human understanding of the items themselves in a natural-language sense. Instead, for each response sequence x , we compute two quantities: $\kappa(x)$ and $\nu(x)$, referred to as the complexity and the surprise parameters respectively (See Methods for precise definitions). Kolmogorov complexity is not computable³², and thus we estimate a related computable parameter (See section on `VeRITAS` analysis in Methods).

The distributions for κ, ν, μ are characterized from the part of the VA data used for training, despite the fact

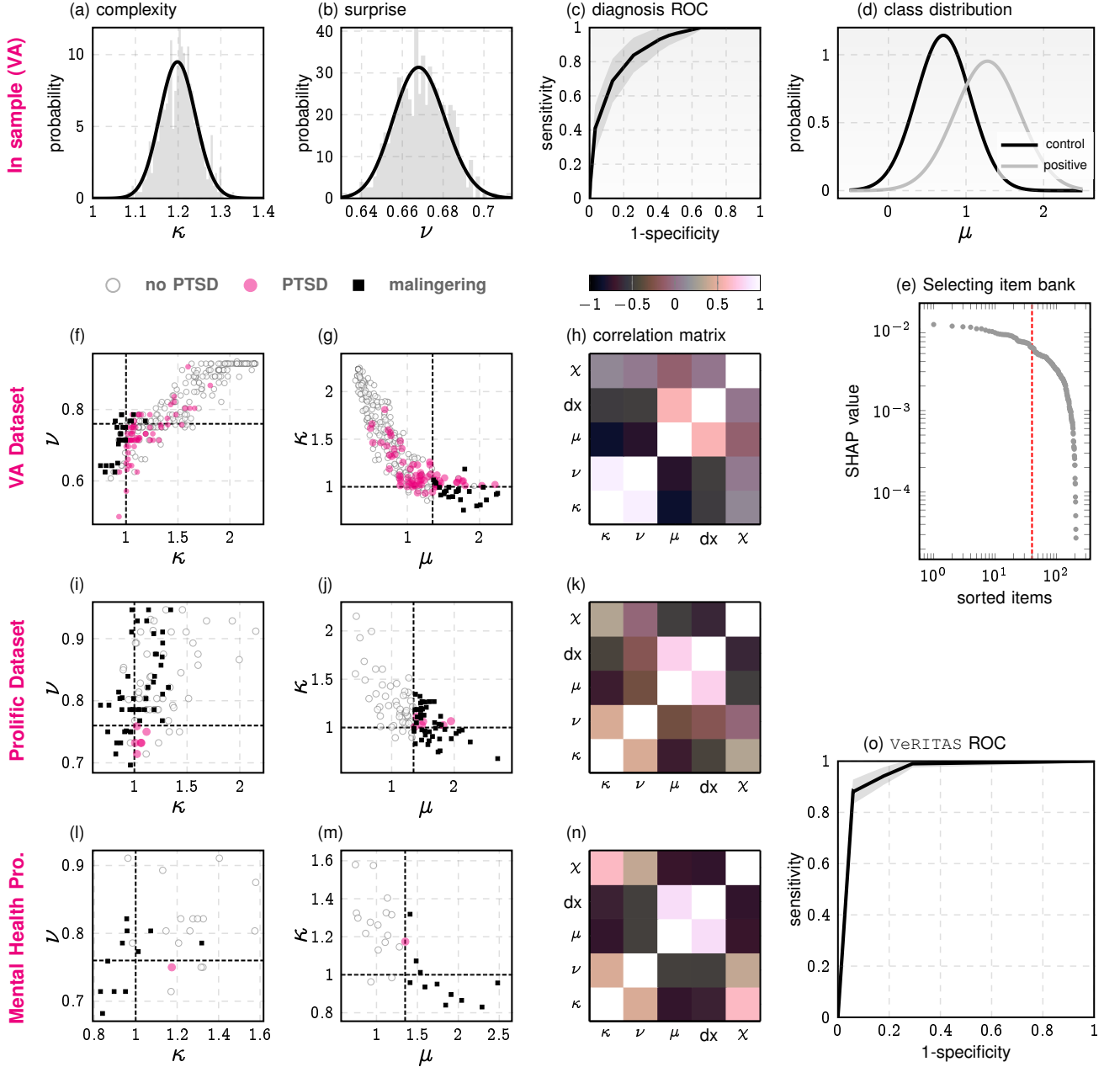


Fig. 2: Training and performance. a, distribution of the complexity parameter. b, distribution of the surprise parameter. c, ROC curve for diagnosis of PTSD without consideration of possible malingering using Q-net models, *i.e.* using the μ parameter as risk. d, Class distribution for setting $\mu_0 = 1.35$. e, distribution of SHAP values of the CAT-PTSD master item bank, showing the threshold above which the VerITAS subset is selected. Panels f-n, out-of-sample results for the three datasets, namely out-of-sample portion of the VA dataset, the PL and the PS datasets. Notably, the different categories of responses obtained comprise ones that are deemed to not have PTSD without any need to consider malingering, those which are diagnosed as having PTSD, and those who are flagged as malingering. Panels h,k,n show correlation between the three VerITAS parameters, and χ which is an indicator variable for malingering flags. dx is the indicator of physician-confirmed diagnosis (only available for the VA data), or estimated diagnosis using the μ_0 threshold. Panel o illustrates the lower envelop of the ROC curve for VerITAS, and 95% confidence bounds.

that this dataset does not have designation of “malingering”. This is possible, because these quantities are computable from just the response sequences themselves. Concretely, we propose a response sequence x should be flagged as an instance of malingering if for suitably chosen thresholds μ_0, ν_0, κ_0 :

$$\chi(x) \triangleq (\mu(x) \geq \mu_0) \bigwedge \left((\kappa(x) \leq \kappa_0) \vee (\nu(x) \geq \nu_0) \right) \quad (1)$$

This may be paraphrased as “a response sequence has high likelihood of malingering if it 1) produces a diagnosis of PTSD with high probability, and 2) is either too surprising or too simple.” The decision thresholds are obtained from theoretical considerations and the VA data (See Methods), which allows us to choose thresholds as reflecting specificity-sensitivity trade-offs (Fig. 2 a,b).

Integrated Diagnostic Capability

Given a sequence of responses to a diagnostic interview, our first task is to determine if a particular subject should be diagnosed with PTSD, if the possibility of malingering is ignored. We call this a “naive diagnosis”. This diagnostic information might be available to `VeRITAS` externally (*e.g.* from a physician’s assessment of the patients). However, `VeRITAS` also has an integrated capability for naive diagnosis: we identify separate generative models for 1) the diagnosed set of patients in the training set (M^+ , some of them might be malingering) and 2) for the patients identified as not having PTSD (M^0). Then given the sequence of responses from a new subject, we estimate if that sequence is more likely to have been generated by the model for M^+ vs that for M^0 . We validate this diagnostic capability using the training data described above, and we achieve good disambiguation between M^+ and M^0 , with out-of-sample AUC 0.867 ± 0.008 (See Fig. 2 c), which is at par or better compared to reported tools, *e.g.*, CAT-PTSD³³. In the `VeRITAS` algorithm, for a given response sequence x , the naive diagnosis risk score is denoted as $\mu(x)$, and referred to as simply the “score” (See Methods). Fig. 2,d shows the class specific distributions estimated for μ for the VA data.

VeRITAS Validation Strategy and Performance

Since the VA data does not indicate presence or absence of malingering, we adopted a non-standard approach for validating `VeRITAS`. We assume that the PL participants do not have PTSD (based on their screening of not having past or present mental health diagnosis, little or no anxiety severity, and being informed to not take the test if experiencing PTSD symptoms), and that all of them were attempting to malingering as directed. This allows us to measure the false negative rate (or 1 - sensitivity), as function of the `VeRITAS` parameters. The success rates of these participants in getting a “diagnosis” at a average sensitivity of **94.2%** (using `VeRITAS` parameters $\kappa_0 = 1, \nu_0 = 0.76, \mu_0 = 1.35$) is shown in Table 1. Then we check how many of the VA participants are flagged as malingering amongst the ones with a PTSD diagnosis. To determine a lower bound on performance we can assume all of these subjects are false positives (which is unlikely, but nevertheless is an upper bound on false positives as function of `VeRITAS` parameters). This allows us to construct a lower envelop of the ROC curve, and hence estimate a lower bound of the AUC for determining malingering, establishing a minimum AUC of **0.95 ± 0.02 at 95% confidence**. High performance operating points for different values of the population prevalence (prevalence has been reported to be high between 20 to 30%)³⁴, reflecting specificity/sensitivity trade-offs are shown in Table 2. These results indicate that if the population prevalence in 30%, then we can achieve **$95.6 \pm 1\%$ sensitivity with 88% specificity, along with $89.4 \pm 0.1\%$ PPV and $97.9 \pm 0.4\%$ NPV, and positive likelihood ratio 9.9 ± 0.9 . For lower population prevalences, *e.g.* at 20%, we can achieve $94.2 \pm 1.2\%$ sensitivity with 93% specificity, along with $83.1 \pm 0.2\%$ PPV and $98.4 \pm 0.3\%$ NPV, and positive likelihood ratio 19.77 ± 0.35 . We note that these numbers represent lower bounds on `VeRITAS` performance, due to our assuming an upper bound on false positives. The variation of the complexity and surprise parameters for the VA data, along with a set of decision thresholds, is shown in Fig. 2 f,g for the VA data, panels i,j for the PL data, and panels l,m for the PS data. The correlation matrix between $\kappa, \nu, \mu, dx, \chi$ (where dx is the variable for clinician diagnosis) for the VA dataset (panel h) shows that the complexity and the surprise parameters are well-correlated, and the dx and the μ parameters are well-correlated, and the malingering decision χ is not very well correlated to either of these. This pattern is closely replicated in the PL and the PS datasets (panels k,n respectively), where the “ dx ” variable refers to the predicted diagnosis. The estimated lower envelop for the ROC curve for the `VeRITAS` algorithm is shown in panel o.**

Item Presentation and Response Sequence in Administration

In `VeRITAS` we present a small random subset of the full item bank to the subjects. Although the models for cross-dependency are inferred using the complete response sequence to all 211 items available, presenting a smaller item set makes the test feasible, while also producing a vast number of possible variations. In the version we validated, 20 items are presented. These are randomly generated from a distribution reflecting the impact of the items in the inferred models on their ability to disambiguate between responses from diagnosed PTSD patients and the control group. **Impact is measured using standard SHAP analysis (See Methods)**. SHAP-values for items are normalized to a probability mass function (pmf) over the top $r = 20$ items, which is then sampled to produce the item set. We can select items adaptively similar to CAT-PTSD³⁰; such complications have been presently avoided. Fig. 2e shows the sorted estimated SHAP values for the items, and shows the threshold of selection for presentation.

Q-net Inference and Cross-talk Modeling

The generative models for response sequences inferred separately from the PTSD-positive and PTSD-negative patients are referred to as “Q-nets”. Structurally, an Q-net comprises an interdependent collection of local predictors, each aiming to predict the response to a specific item using as features the responses to other items from the item bank. (Fig. 1). Thus, an Q-net comprises almost as many such position-specific predictors as the length of the response sequence. These individual predictors are implemented as conditional inference

trees³⁵, in which nodal splits have a minimum pre-specified significance in differentiating the child nodes. Thus, each predictor yields an estimated conditional response distribution for each item. The set of items acting as features in each predictor are automatically identified, *e.g.*, in the fragment of PTSD-positive Q-net (Fig 1), the predictor for item `ptsd119` includes the response to item `ptsd19`, that for `ptsd19` includes item `ptsd43`, that for `ptsd43` includes `ptsd208`, which in turn includes `ptsd43`, revealing complex possibly cyclic dependencies. Owing to this recursive expansion, a complete Q-net captures the complexity of the rules guiding the subtle cross-talk emerging between different items in a survey, and can be trained without linguistic understanding of the items themselves. The Q-nets inferred for M^+ and M^0 have 5,852 and 10,301 independent parameters, encoding the probability distributions on the tree leaves, which were identified automatically in the course of inferring the sets of conditional inference trees. The overall `VeRITAS` model, which comprises these Q-nets along with the three thresholds described before, therefore has 16,156 parameters.

Time-to-completion and Malingering

In the PL dataset, the average completion time was 3.4 minutes ([2.9, 3.9], 95% CI). We find that malingering flags have no significant correlation with the time to completion. The success rate of defeating `VeRITAS` is observed to have a correlation of 0.48([-.27, .87], 95% CI). Thus, fast responses that might reflect poor quality survey in general^{36–38} might not be effective in detecting malingering, especially if participants are putting more thought into how to fake their symptoms.

Performance among Mental Health Professionals

The PS dataset illustrates that training in forensic psychiatry or neuropsychology confers no significant advantage in defeating the `VeRITAS` algorithm. Only one participant succeeded in not being flagged, implying a false negative rate of 3.7%. Using the same settings, participants in the PL dataset (drawn from the general population in the US and the UK as part of Prolific’s cohort), were able to achieve a false negative rate of 5.8%, which is about 56.8% worse compared to the cohort comprising a substantial number of participants with training in forensic psychiatry.

Number of Distinct Variations of `VeRITAS` Implementation

With $r = 20$ items presented, randomly chosen from the top $N = 40$ shortlisted from the master item bank (sorted according to SHAP values, as described before), and each item having $L = 5$ possible responses, we can have $\binom{N}{r} = 137,846,528,820$ variations of the test, with $L^r \binom{N}{r} = 13,146,069,414,138,793,945,312,500$ or approximately 10^{25} possible responses, which is approximately equal to the number of stars in the observable universe. Thus, it is non-trivial for human subjects to “learn” or “train” to defeat the algorithm.

Comparison Against State of Art

A well-known tool used to identify malingering in mental health diagnoses is the Structured Interview of Reported Symptoms, 2nd Edition (SIRS-2)²³, which has a reported performance of sensitivity of 0.80, a specificity of 0.975, and positive and negative predictive powers of more than 0.90 (based on a base rate of 31.8%), takes 30-40 minutes to complete, needs extensive expert interpretation, and is not disease-specific. In contrast, `VeRITAS` may be completed in under 4 minutes, can have sensitivity and specificity both above 90%, has PPV over 86% and NPV over 90% in selected operating points, can be tuned to specific disorders, and may be administered automatically. The crucial difference in `VeRITAS` is the near-impossibility of defeating it through coaching, and its effectiveness in the scenario that the subject has training as a mental health professional. The principles on which existing tools such as the SIRS-2 are based makes them highly unlikely to be effective if the subject is familiar with the symptomologies of mental disorders, and the approaches employed to flag malingering.

DISCUSSION

The `VeRITAS` algorithm aims to identify feigned, faked, or simulated symptoms in clinician rated and patient self-reported evaluations. While we demonstrate applicability in assessments for PTSD, the underlying principles are generally applicable for detecting such simulated symptoms for other mental health disorders, and even more generally for vetting possibly adversarial responses in structured interviews.

Our findings are relevant to the detection of feigned, faked, or simulated symptoms. Although these are traditionally thought to reflecting malingering, it is important to understand that the commonly applied construct of malingering has led to a misunderstanding of the nature of simulated symptoms. At a descriptive level, simulated symptoms are characterized by a response style that prevents the accurate measure of symptom severity in the evaluation of a medical syndrome. Simulated symptoms are found in factitious disorder and malingering. Factitious disorder is a psychiatric condition which involves simulated symptoms. But unlike in the case of malingering, feigning poor health is thought to be unintentional insofar as there are no clear positive incentives

for illness and there the person is unaware of the fact that they do not have a medical disorder. Malingering is not a diagnostic category in DSM-5 but rather is coded in a special section on clinical phenomena that are not well understood and/or do not represent a mental illness (DSM-5-TR for reference). While the lay language of malingering suggests that with regards to clinical deception, “you know it when you see it”, empirical work on malingering suggests a far more complex picture than that of a criminally motivated deception. The three subtypes of malingered PTSD symptoms include manufactured, exaggerated, and misattributed³⁹. Malingered symptoms can be motivated by criminological, pathogenic, or adaptational motivations⁴⁰. For example, in the emergency department a patient with a highly problematic substance use disorder may accurately judge that a clinician cannot be trusted to provide an unbiased assessment of their need for care and resources. In such a case, a short term strategy with adaptive value could be to simulate symptoms. In the long term, the strategy would be maladaptive and further erode trust between the patient and clinician. In fact, empirical studies of malingering in the emergency department have shown evidence that far from faking bad and lacking psychopathology, persons who feign symptoms of a psychiatric disorder are at higher risk for psychopathology, mortality, are more likely to be homeless or Black/African-American, and more likely to have a substance use disorder than matched controls^{20,41}. Thus, the problem of simulated symptoms and inaccurate diagnosis is not one of detecting malingering in order to deny resources, but rather the prevention of potentially harmful side effects on inappropriate treatments such as medications or involuntary psychiatric hospitalization, and a likely missed opportunity to reduce the risk of mortality and morbidity by meeting the true needs of the person feigning symptoms. Additionally, the measurement of PTSD symptom severity is used to guide treatment decisions. Invalid responses in a person with true PTSD would lead to potentially harmful treatment escalation.

The `VeRITAS` algorithm introduced here aims to identify deception in clinician rated and patient self-reported symptoms associated. While we demonstrate applicability in assessments for PTSD, the underlying principles are generally applicable for detecting faked symptoms for other mental health disorders, and even more generally for vetting adversarial responses in structured interviews.

Malingering presents a significant challenge in accurate diagnosis of mental health disorders. Estimates of prevalence of faked symptoms for PTSD in psychiatric and criminal justice settings range from 8% to 64%^{34,42,43}. This variability in reported estimates highlights the complexity of the issue and the need for reliable, rigorous, and principled methods for detecting deception in this context⁴⁴.

Detection methods for mental health conditions, and for PTSD in particular, must sensitively navigate the complex motivations behind their potential fabrication. These motivations extend beyond financial incentives such as disability compensation, insurance claims, or legal benefits, encompassing mental health issues and socio-economic factors that restrict access to healthcare. The presence of malingering, particularly when prevalent in certain clinics, may point to broader societal issues. Furthermore, a high incidence of malingering not only burdens healthcare and legal systems but also compromises the integrity of clinical diagnoses and research. This misalignment may result in the misallocation of resources and impede the treatment of those genuinely suffering from PTSD. Identifying such behaviors efficiently and discreetly, avoiding the stigma of a “malingering test”, is crucial in addressing the multifaceted challenges encountered in clinical practice. Research indicates the financial and clinical significance of accurately diagnosing mental health conditions, with disability payments and the prevalence of service-connected mental disorders among veterans highlighting the issue’s complexity. Moreover, the prevalence of malingering varies by context, underscoring the need for a nuanced approach to diagnosis and treatment.

Our approach represents a significant departure from traditional methods of malingering detection, which often rely on domain-specific knowledge⁴⁵ and standardized tests developed through extensive research on both genuine patients and malingerers (Table 4). While these traditional methods, including the SIRS²³ (requiring 30-40 minutes and expert assessments), the Structured Inventory of Malingered Symptomology (SIMS)²⁴, and validity scales associated with the Minnesota Multiphasic Personality Inventory-2²⁵, have shown accuracy rates in the range of 85% to 95% depending on the specific problem context^{46–49}, they come with limitations. These include the requirement for substantial expertise to develop and administer, vulnerabilities to clinician bias, and the potential for false positives and negatives. Furthermore, they may not easily extend to other contexts or disorders, and their effectiveness can be compromised by coaching or prior knowledge of psychiatric symptomology. Additionally, current practice does not have good methods to discreetly flag malingering. Attempting to make such assessments without any formal tools is problematic, since humans are generally poor at detecting lies, with accuracy rates often barely surpassing chance²⁶.

In contrast, `VeRITAS` leverages statistical differences in an individual’s response patterns compared to a set of baseline responses, aiming to capture the complex, often non-obvious dependencies between interview questions. The underlying model is a nonparametric generative model that maps the inter-question relations and dependencies, enabling the detection of inconsistencies indicative of malingering.

Our findings indicate that `VeRITAS` can achieve high sensitivity and specificity in detecting malingering, with performance metrics potentially surpassing or at least comparable to those of existing state-of-the-art techniques, but with several advantages. `VeRITAS` requires less time for administration, does not necessitate domain-

specific expertise for interpretation, and minimizes the risk of bias. Moreover, its design makes it challenging for individuals, even those with psychiatric training, to defeat the system through coaching or preparation.

However, the impact of `VeRITAS` in clinical practice must be evaluated carefully, particularly for unintended ethical implications, especially with respect to vulnerable communities within which the impact of mental health disorders are often exacerbated by limited access to healthcare, socio-economic instability, and the stigma surrounding mental health diagnoses. A non-zero risk of false positives could unjustly exclude vulnerable individuals from receiving necessary care, and the algorithm’s reliance on statistical patterns may overlook the nuanced expressions of PTSD symptoms across different cultures, possibly leading to biased assessments. This underscores the importance of integrating cultural sensitivity into the algorithm’s development and emphasizing the indispensable role of human judgment in the diagnostic process^{26,44}. Ongoing research and validation in diverse settings are imperative to refine the algorithm’s application, ensuring it serves as a tool for empowerment rather than exclusion, particularly in underserved populations.

In conclusion, while malingering presents a persistent challenge in the accurate diagnosis and treatment of PTSD, our study offers a promising new direction for detection methods. By employing a sophisticated, data-driven approach that transcends the limitations of traditional methods, `VeRITAS` provides a powerful tool for clinicians and researchers. Its adoption could significantly enhance the integrity of PTSD diagnoses, ensuring that resources are allocated to those genuinely in need and supporting the broader goals of psychiatric care and research. However, further research and validation across diverse populations and settings are essential to fully realize its potential and applicability in clinical practice.

METHODS

1. DEFINITIONS & NOTATION

Definition 1 (Survey). *A survey for the purpose of this work is a structured interview, consisting of a finite number of questions (items) posed to a set of participants, with these items drawn from a finite item bank, and whose responses must be one from a pre-specified set of choices, e.g., the Likert scale, with missing values for the responses allowed.*

Definition 2 (Response vector). *A response vector is the set of responses to a survey from a single participant, typically assuming that not all items are posed, and allows for the possibility that some responses are missing.*

A Q-net, as described here, is a model of the response dependency structure for questions (items) posed to participants in a survey. The Q-net explicitly estimates individual conditional distributions of each item response, which collectively serve as a model of the full joint distribution of the responses.

Definition 3 (Q-net). *Let $X \sim P$ be an n -dimensional discrete random vector supported on a finite set Σ and following distribution P , i.e.*

$$X = (X_1, \dots, X_n) \sim P, \quad \text{supp}(X) = \Sigma = \prod_{i=1}^n \Sigma_i \quad \text{with } |\Sigma| < \infty.$$

For $i = 1, \dots, n$, let $P_i := P(X_i | X_j = x_j \text{ for } j \neq i)$ denote the conditional distribution of X_i given the values of the other components of X . Finally, for each $i = 1, \dots, n$, let Φ_i^P denote an estimate of the distribution P_i . Then the set $\Phi^P := \{\Phi_i^P\}_{i=1}^n$ is called a Quasinet (Q-net). Identifying the true distribution P as the one describing the joint statistics of the responses from a survey with n items, we also refer to Φ^P as the Q-net for the survey P .

When P is clear from context, we may omit the superscript and simply write $\Phi = \{\Phi_i\}$ to denote the Q-net. The motivation for Definition 3 is that the collection of all estimators $\Phi = \{\Phi_i\}$ contained in a Q-net represents the set of all inferred dependencies from the observed ecosystem. While the definition allows for arbitrary method of algorithm to construct the estimators Φ_i , the utility of a Q-net clearly depends primarily on the properties of the Φ_i . In this study, we aim to minimize the set of a priori assumptions on the overall model structure to allow the complex dependencies present in P to emerge. To that end, throughout this work all Q-nets are computed using conditional inference trees⁵⁰ (a variant of classification and regression trees) to compute each Φ_i . In general each Q-net component Φ_i is computed independently from the other Φ_j , which allows a network structure to emerge amongst these estimators.

An important quantity for an inferred Q-net is the persistence function ω_x .

Definition 4 (Persistence Function). *Given a survey P inducing the Q-net Φ^P and a response vector $x = (x_1, \dots, x_n)$, the persistence ω_x of x in the population modeled by the Q-net:*

$$\omega_x^P := \Pr(x \in P) = \prod_{i=1}^n \Phi_i^P(X_i = x_i | X_j = x_j, j \neq i) \quad (2)$$

The persistence function ω_x^P , as the name suggests, is the probability that x persists, i.e., $Pr(x \rightarrow x)$ for the population modeled by the Q-net P , with $1 - \omega_x^P$ being the probability that x is altered by a random perturbation.

We will show that if for two inferred Q-net models P, Q , we have $\omega_x^P \geq \omega_x^Q$, then it is more likely that model P generated x . This is an important result that justifies the definition of the score parameter in Defn. 6.

The Q-net allows us to rigorously compute bounds on the probability of a spontaneous change from one response vector to another, induced by spontaneous chance variations. Not all perturbations in a vector are either likely or contextually meaningful. With an exponentially exploding number of possibilities in which a vector over a large set of items can vary, it is computationally intractable to directly model all possible dependencies; nevertheless, we can constrain the possibilities using the patterns we uncover via the Q-net construction. A key piece of this approach is to design an intrinsic distance (q-distance) between any two response vectors, which is reflective of this underlying dependency structure.

Definition 5 (q-distance). Let $\Phi^P = \{\Phi_i^P\}_{i=1}^n$ and $\Phi^Q = \{\Phi_i^Q\}_{i=1}^n$ denote Q-nets on populations P and Q , and suppose $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$ are samples of $X \sim P$ and $Y \sim Q$ respectively. Then the q-distance $\theta_{P,Q}(x, y)$ between x and y is

$$\theta_{P,Q}(x, y) := \frac{1}{n} \sum_{i=1}^n \left[\mathbb{J}^{\frac{1}{2}} \left(\Phi_i^P(X_i | X_j = x_j, j \neq i) \parallel \Phi_i^Q(Y_i | Y_j = y_j, j \neq i) \right) \right]$$

where \mathbb{J} denotes the Jensen-Shannon divergence⁵¹.

For brevity, we may write simply θ (dropping the suffixes) if the populations are clear from context. Since the Jensen-Shannon distance \mathbb{J} is a legitimate metric⁵² on the set of probability distributions (unlike KL-divergence), θ inherits nonnegativity, symmetry, and respects the triangle inequality; it follows that q-distance is a (pseudo)-metric on Σ . Note that, being a pseudo-metric implies that we may have $\theta(x, y) = 0$ for $x \neq y$, i.e. distinct vectors can induce the same distributions over each index, and thus have zero distance. This is in fact desirable, since we do not want our distance to be sensitive to changes that are not meaningful. The intuition is that not all variations are equally important or likely. Moreover, we show in Theorem 1 that the log-likelihood of a vector x transitioning to y scales with $\theta(x, y)$, allowing us to directly estimate the probability of spontaneous (or sequential) jumps between abundance profiles.

Theorem 1 (Probability Bound). Given a vector x of length n from P that transitions to y from Q , we have the following bounds at significance level α .

$$\omega_y e^{\frac{\sqrt{8N^2}}{1-\alpha} \theta(x, y)} \geq Pr(x \rightarrow y) \geq \omega_y e^{-\frac{\sqrt{8N^2}}{1-\alpha} \theta(x, y)} \quad (3)$$

where ω_y is the persistence of y (Def. 4), and $\theta(x, y)$ is the q-distance between x, y (Def. 5).

Proof. See later in Section 3. □

Theorem 1 gives theoretical backing to the claim that samples generated by the Q-net indeed reflect likely perturbation possibilities from the current state. Thus we can use the Q-net to draw contextually realistic samples that respect the cross dependencies and reduce surprise (that is, the Q-net-inferred conditional distributions can be used to generate approximate samples from the population P).

Remark 1 (Neighborhood Structure). It follows from Th. 1 that we have for some constant C ,

$$\ln \left| \frac{Pr(x \rightarrow y)}{Pr(y \rightarrow x)} \right| \leq C \theta(x, y) \quad (4)$$

implying for all response vectors y within a small neighborhood of x (small in metric θ), we have:

$$Pr(y \rightarrow x) \approx Pr(x \rightarrow x) \quad (5)$$

which reveals an important special structure on local neighborhoods.

2. VERITAS ANALYSIS

Definition 6 (Algorithm VERITAS Parameters). We introduce three parameters referred to as the complexity, surprise and score parameters (κ, ν, μ respectively) for a given response vector x :

$$\text{complexity: } \kappa \triangleq -\frac{1}{|x|} \ln Pr(x \rightarrow x | M^+) = -\frac{\ln \omega_x^{M^+}}{|x|} \quad (6)$$

$$\text{surprise: } \nu \triangleq \mathbf{E}_i \left(1 - \Phi_i^{M^+}(x_{-i}) \mid x_i \right) \quad (7)$$

$$\text{score: } \mu \triangleq \frac{\ln Pr(x \rightarrow x | M^+)}{\ln Pr(x \rightarrow x | M^0)} = \frac{\ln \omega_x^{M^+}}{\ln \omega_x^{M^0}} \quad (8)$$

where M^+ indicates the sub-population exhibiting a particular trait of interest e.g. a mental health disorder such as PTSD, and M^0 is the control sub-population where this trait is absent.

Definition 7 (Malingering property). A response vector x is defined to have the malingering property if:

$$\chi(x) \triangleq (\mu(x) \geq \mu_0) \bigwedge \left((\kappa(x) \leq \kappa_0) \vee (\nu(x) \geq \nu_0) \right) \quad (9)$$

The decision thresholds κ_0, ν_0, μ_0 are inferred from survey data.

Lemma 1 (Complexity). For a survey with n items, and assuming L to be the number of possible responses to each item, the unconditional probability of a response vector x occurring among all feasible responses is bounded above by $(e^\kappa/L)^n$, where $\kappa(x)$ is the complexity parameter for response x .

Proof. Let $\kappa(x) \leq \kappa'$. From Def. 6, we have for a response vector x ,

$$-\frac{1}{n} \ln \omega_x \leq \kappa' \Rightarrow \omega_x \geq e^{-n\kappa'} \quad (10)$$

Summing on both sides over all responses x with $\kappa(x) \leq \kappa'$ (assume there are N_x such sequences), we have:

$$1 \geq \sum_x \omega_x \geq \sum_x e^{-n\kappa'} \quad (11)$$

where the first inequality follows from observing that responses very close to x in the q-distance metric have a specific structure, namely $\omega_x \approx \Pr(y \rightarrow x)$ (See Remark 1) and responses further away have smaller jump probabilities, which then implies:

$$N_x \sum_x e^{-n\kappa'} \leq 1 \Rightarrow N_x \leq e^{n\kappa'} \quad (12)$$

The result then follows from noting that the complete set of possible responses has the size L^n . \square

Lemma 1 justifies why a low value of κ implies the possibility of an un-natural response, because the odds of generating such a response is remarkably small.

Corollary 1 (Algorithmic Complexity). , The algorithmic complexity of a response x conditional on the number of survey items n is at most $\kappa(x) + O(1)$.

Proof. This follows from noting that a set of cardinality L^m has a algorithmic complexity of $m + O(1)$, since words of length m are sufficient to encode the index of any element of the set, and thus can be uniquely identified. Sinec we can calculate $\kappa' = \kappa(x)$ for any x , and since the set of all x for a given value of κ' belongs to a set of size at most $e^{-n\kappa'}$, the result follows. \square

Lemma 2 (Surprise). For any response vector x , we have:

$$\nu(x) \leq 1 - e^{-\kappa(x)} \quad (13)$$

Proof. Denoting $\Phi_i(x_{-i})|_{x_i}$ as a_i , we note that $\omega_x^{1/n}$ is the geometric mean of the vector of a_i s, while $\mathbf{E}_i(\Phi_i(x_{-i})|_{x_i})$ is the arithmetic mean of the same vector, which then completes the proof by noting:

$$-\mathbf{E}_i(\Phi_i(x_{-i})|_{x_i}) \leq -\omega_x^{1/n} \Rightarrow \nu(x) \leq 1 - \omega_x^{1/n} \quad (14)$$

\square

Interpretation on Why the Defined Property Identifies Malingering

Lemma 2 indicates that the requirement of an upper bound on the surprise and a lower bound on the complexity are both aiming to flag responses which are unlikely to appear when the data (responses) are being generated by the underlying process corresponding to the phenotype of interest (PTSD). When such unlikely responses do appear nevertheless, it is likely that they are not being generated by the correct underlying process. One can attempt to fake responses that might seem to increase the odds of a positive diagnosis, but the respondent must replicate the cross-dependencies closely enough (build in enough structure) so that the deviation from the expected responses is limited (limited surprise requirement). But building in too much structure will reduce the complexity too much (too much structure reduces complexity, since there are fewer highly structured sequences), which will then fail the complexity lower bound.

Note that the remaining condition $\mu(x) \geq \mu_0$ is a diagnosis criterion for the trait of interest (M^+), and may be replaced with a different condition if available for identifying participants with the M^+ trait. This particular form follows from a straightforward Bayesian argument on estimating the posterior.

SHAP Analysis Selection of Item Bank

We present a random subset of $r = 20$ items from a master subset of items used in CAT-PTSD. The items that are selected to make this master subset is obtained by ranking the items according to their estimated impact on model prediction. One standard approach to estimate impact of features on model outcomes is SHAP (SHapley Additive exPlanations) analysis⁵³, which is a method derived from game theory to explain the output of machine learning models. It provides a way to measure the contribution of each feature in a given model to the prediction for each instance. In our scenario, we use the persistence function as the model prediction to compute SHAP values, *i.e.*, we rank the items based on the degree to which including an item within a subset of items with non-empty responses moves the value of the persistence function.

3. PROOF OF THEOREM 1

Theorem 2 (Probability bound). *Given a sequence x of length N that transitions to a strain $y \in Q$, we have the following bounds at significance level α .*

$$\omega_y^Q e^{\frac{\sqrt{8}N^2}{1-\alpha}\theta(x,y)} \geq Pr(x \rightarrow y) \geq \omega_y^Q e^{-\frac{\sqrt{8}N^2}{1-\alpha}\theta(x,y)} \quad (15)$$

where ω_y^Q is the membership probability of strain y in the target population Q (See Def. ??), and $\theta(x, y)$ is the q -distance between x, y (See Def. 5).

Proof. Using Sanov's theorem⁵¹ on large deviations, we conclude that the probability of spontaneous jump from strain $x \in P$ to strain $y \in Q$, with the possibility $P \neq Q$, is given by:

$$Pr(x \rightarrow y) = \prod_{i=1}^N (\Phi_i^P(x_{-i})|_{y_i}) \quad (16)$$

Writing the factors on the right hand side as:

$$\Phi_i^P(x_{-i})|_{y_i} = \Phi_i^Q(y_{-i})|_{y_i} \left(\frac{\Phi_i^P(x_{-i})|_{y_i}}{\Phi_i^Q(y_{-i})|_{y_i}} \right) \quad (17)$$

we note that $\Phi_i^P(x_{-i})$, $\Phi_i^Q(y_{-i})$ are distributions on the same index i , and hence:

$$|\Phi_i^P(x_{-i})|_{y_i} - \Phi_i^Q(y_{-i})|_{y_i}| \leq \sum_{y_i \in \Sigma_i} |\Phi_i^P(x_{-i})|_{y_i} - \Phi_i^Q(y_{-i})|_{y_i}| \quad (18)$$

Using a standard refinement of Pinsker's inequality⁵⁴, and the relationship of Jensen-Shannon divergence with total variation, we get:

$$\theta_i \geq \frac{1}{8} |\Phi_i^P(x_{-i})|_{y_i} - \Phi_i^Q(y_{-i})|_{y_i}|^2 \Rightarrow \left| 1 - \frac{\Phi_i^Q(y_{-i})|_{y_i}}{\Phi_i^P(x_{-i})|_{y_i}} \right| \leq \frac{1}{a_0} \sqrt{8\theta_i} \quad (19)$$

where a_0 is the smallest non-zero probability value of generating the entry at any index. We will see that this parameter is related to statistical significance of our bounds. First, we can formulate a lower bound as follows:

$$\log \left(\prod_{i=1}^N \frac{\Phi_i^P(x_{-i})|_{y_i}}{\Phi_i^Q(y_{-i})|_{y_i}} \right) = \sum_i \log \left(\frac{\Phi_i^P(x_{-i})|_{y_i}}{\Phi_i^Q(y_{-i})|_{y_i}} \right) \geq \sum_i \left(1 - \frac{\Phi_i^Q(y_{-i})|_{y_i}}{\Phi_i^P(x_{-i})|_{y_i}} \right) \geq \frac{\sqrt{8}}{a_0} \sum_i \theta_i^{1/2} = -\frac{\sqrt{8}N}{a_0} \theta \quad (20)$$

Similarly, the upper bound may be derived as:

$$\log \left(\prod_{i=1}^N \frac{\Phi_i^P(x_{-i})|_{y_i}}{\Phi_i^Q(y_{-i})|_{y_i}} \right) = \sum_i \log \left(\frac{\Phi_i^P(x_{-i})|_{y_i}}{\Phi_i^Q(y_{-i})|_{y_i}} \right) \leq \sum_i \left(\frac{\Phi_i^Q(y_{-i})|_{y_i}}{\Phi_i^P(x_{-i})|_{y_i}} - 1 \right) \leq \frac{\sqrt{8}N}{a_0} \theta \quad (21)$$

Combining Eqs. 20 and 21, we conclude:

$$\omega_y^Q e^{\frac{\sqrt{8}N}{a_0}\theta} \geq Pr(x \rightarrow y) \geq \omega_y^Q e^{-\frac{\sqrt{8}N}{a_0}\theta} \quad (22)$$

Now, interpreting a_0 as the probability of generating an unlikely event below our desired threshold (*i.e.* a "failure"), we note that the probability of generating at least one such event is given by $1 - (1 - a_0)^N$. Hence if α is the pre-specified significance level, we have for $N \gg 1$:

$$a_0 \approx (1 - \alpha)/N \quad (23)$$

Hence, we conclude, that at significance level $\geq \alpha$, we have the bounds:

$$\omega_y^Q e^{\frac{\sqrt{8}N^2}{1-\alpha}\theta} \geq Pr(x \rightarrow y) \geq \omega_y^Q e^{-\frac{\sqrt{8}N^2}{1-\alpha}\theta} \quad (24)$$

□

TABLE 4: Summary of malingering/deception detection methods.

Approach/Tool	Description	Noted Accuracy/Performance
---------------	-------------	----------------------------

Structures Interview of Reported Symptoms (SIRS) ²³	An interview-based measure with multiple detection strategies.	No specific accuracy rate mentioned, but noted as a robust instrument.
Structured Inventory of Malingered Symptomology (SIMS) ²⁴	A paper-and-pencil screening device for detecting malingering with a sensitivity for detecting malingering of 95.6% in a study with college students.	95.6% sensitivity in a specific study context.
Minnesota Multiphasic Personality Inventory-2 (MMPI-2) ²⁵	A self-report measure assessing personality and psychopathology. Certain validity scales were developed to uncover malingering.	Not specified, but noted flaws and potential for false positives.
Millon Clinical Multiaxial Inventory MCMI-III	A self-report scale focusing on personality disorders.	No specific accuracy rate mentioned.
Miller Forensic Assessment of Symptoms (M-FAST)	A brief screening measure for malingered mental illness in forensic settings.	Noted issues with low internal consistency on some scales.
Human Lie Detectors	Not a specific tool, but the general method of humans attempting to discern lies from truth.	People are generally poor lie detectors.
Arousal-Based Approaches (like the Polygraph)	Techniques that rely on physiological responses.	Criticized for poor validity and high rate of false positives.
Cognitive Load-Inducing Approaches	Techniques that view deception as a cognitive act that generally imposes greater cognitive load on respondents than honesty does.	No specific accuracy rate mentioned.
Autobiographical Implicit Association Test (aIAT) ⁴⁷	Designed to determine whether respondents possess actual autobiographical memories.	91% accuracy rate in identifying genuine autobiographical memories.
Timed Antagonistic Response Alethiometer (TARA) ⁴⁸	A computer-administered, response time-based method of lie detection.	85% accuracy rate.
Detecting Faked Identities With Unexpected Questions and Mouse Movements	Technique using computer mouse movements in conjunction with unexpected questions to uncover faked identities ⁴⁹	95% accuracy rate.
Time-Restricted Integrity Confirmation (TRI-Con)	A cognitive load-inducing technique with potential to uncover different kinds of deception including malingering.	Up to 89% accuracy rate.
Activation-Decision-Construction-Action Theory (ADCAT)	A theory of high-stakes deception.	No specific accuracy rate mentioned, but this is more of a theoretical foundation rather than a specific tool or method.

REFERENCES

- [1] Rogers, R. An overview of malingering and its assessment. *Psychiatric Clinics of North America* **20**, 15–27 (1997).
- [2] Rogers, R. E. *Clinical assessment of malingering and deception* (Guilford Press, 2008).
- [3] Frueh, B., Grubaugh, A., Elhai, J. & Buckley, T. Us department of veterans affairs disability policies for posttraumatic stress disorder: administrative trends and implications for treatment, rehabilitation, and research. *Am J Public Health* **97**, 2143–5 (2007).
- [4] Taylor, S., Frueh, B. C. & Asmundson, G. J. G. Detection and management of malingering in people presenting for treatment of posttraumatic stress disorder: Methods, obstacles, and recommendations. *Journal of Anxiety Disorders* **21**, 22–41 (2007).

- [5] Ali, S., Jabeen, S. & Alam, F. Multimodal approach to identifying malingered posttraumatic stress disorder: A review. *Innovations in Clinical Neuroscience* **12**, 12 (2015).
- [6] Butcher, J. N. Minnesota multiphasic personality inventory. *The Corsini Encyclopedia of Psychology* 1–3 (2010).
- [7] Gurriel, J. & Fremouw, W. Assessing malingered posttraumatic stress disorder: A critical review. *Clinical Psychology Review* **23**, 881–904 (2003).
- [8] Salloway, S., Southwick, S. & Sadowsky, M. Opiate withdrawal presenting as posttraumatic stress disorder. *Hospital and Community Psychiatry* **41**, 666–667 (1990).
- [9] Resnick, P. J., West, S. & Payne, J. W. Malingering of posttraumatic disorders. In Rogers, R. (ed.) *Clinical assessment of malingering and deception*, 109–127 (Guilford Press, 2008), 3 edn.
- [10] Burkett, B. G. & Whitley, G. *Stolen valor: How the Vietnam generation was robbed of its heroes and history* (Verity Press, 1998).
- [11] Goldstein, R. B. *et al.* The epidemiology of dsm-5 posttraumatic stress disorder in the united states: results from the national epidemiologic survey on alcohol and related conditions-iii. *Social psychiatry and psychiatric epidemiology* **51**, 1137–1148 (2016).
- [12] Schnurr, P. P., Lunney, C. A., Bovin, M. J. & Marx, B. P. Posttraumatic stress disorder and quality of life: Extension of findings to veterans of the wars in iraq and afghanistan. *Clinical psychology review* **29**, 727–735 (2009).
- [13] LoPiccolo, C., Goodkin, K. & Baldewicz, T. Current issues in the diagnosis and management of malingering. *Ann Med* **31**, 166–174 (1999).
- [14] Oboler, S. Disability evaluations under the department of veterans affairs. In Rondinelli, R. & Katz, R. (eds.) *Impairment Rating and Disability Evaluation*, 187–217 (W. B. Saunders, Philadelphia, PA, 2000).
- [15] Taylor, S. *Clinician's Guide to Treating PTSD: A Cognitive-Behavioral Approach* (Guilford Press, New York, 2006).
- [16] Rosen, G. Dsm's cautionary guideline to rule out malingering can protect the ptsd data base. *J Anxiety Disorders* **20**, 530–535 (2006).
- [17] Marx, B. & Holowka, D. Ptsd disability assessment. *PTSD Res Q* **22**, 1–6 (2011).
- [18] Rogers, R., Sewell, K. & Goldstein, A. Explanatory models of malingering: a prototypical analysis. *Law & Hum Behav* **18**, 543–52 (1994).
- [19] Lees-Haley, P. Mmpi-2 base rates for 492 personal injury plaintiffs: implications and challenges for forensic assessment. *J Clin Psychol* **53**, 745–55 (1997).
- [20] Park, L., Costello, S., Li, J., Lee, R. & Jacobson, K. C. Race, health, and socioeconomic disparities associated with malingering in psychiatric patients at an urban emergency department. *General Hospital Psychiatry* **71**, 121–127 (2021).
- [21] Muntaner, C., Eaton, W. W., Miech, R. & O'campo, P. Socioeconomic position and major mental disorders. *Epidemiologic reviews* **26**, 53–62 (2004).
- [22] Drob, S. L., Meehan, K. B. & Waxman, S. E. Clinical and conceptual problems in the attribution of malingering in forensic evaluations. *The journal of the American Academy of Psychiatry and the Law* **37**, 98–106 (2009).
- [23] Wong, S. & O'Sullivan, M. The structured interview of reported symptoms (sirs): An overview. *Assessment* **12**, 289–307 (2005).
- [24] Smith, G. P. & Burger, G. K. Detection of malingering: validation of the structured inventory of malingered symptomatology (sims). *Journal of the American Academy of Psychiatry and the Law Online* **25**, 183–189 (1997).
- [25] Ben-Porath, Y. S. *Interpreting the mmpi-2-rf* (U of Minnesota Press, 2012).
- [26] Ekman, P. & O'Sullivan, M. Who can catch a liar? *American Psychologist* **46**, 913 (1991).
- [27] Mihalcea, R. & Strapparava, C. The lie detector: Explorations in the automatic recognition of deceptive language. In *ACL-IJCNLP 2009* (2009).
- [28] Burgoon, J. K., Blair, J. P. & Strom, R. E. Cognitive biases and nonverbal cue availability in detecting deception. *Human Communication Research* **34**, 572–599 (2008).
- [29] Zhou, L., Burgoon, J. K., Nunamaker Jr, J. F. & Twitchell, D. Automating linguistics-based cues for detecting deception in text-based asynchronous computer-mediated communications. *Group Decision and Negotiation* **13**, 81–106 (2004).
- [30] Brenner, L. A. *et al.* Development and validation of computerized adaptive assessment tools for the measurement of posttraumatic stress disorder among us military veterans. *JAMA Network Open* **4**, e2115707–e2115707 (2021).
- [31] Gibbons, R. D. *et al.* Development of a computerized adaptive test for depression. *Archives of general psychiatry* **69**, 1104–1112 (2012).
- [32] Li, M., Vitányi, P. *et al.* *An introduction to Kolmogorov complexity and its applications*, vol. 3 (Springer, 2008).
- [33] Brenner, L. A. *et al.* Development and validation of computerized adaptive assessment tools for the measurement of posttraumatic stress disorder among us military veterans. *JAMA Network Open* **4**, e2115707 (2021). URL <https://doi.org/10.1001/jamanetworkopen.2021.15707>.

-
- [34] Matto, M., McNeil, D. E. & Binder, R. L. A systematic approach to the detection of false ptsd. *The journal of the American Academy of Psychiatry and the Law* **47**, 325–334 (2019).
 - [35] Hothorn, T., Hornik, K. & Zeileis, A. Unbiased recursive partitioning: A conditional inference framework. *JOURNAL OF COMPUTATIONAL AND GRAPHICAL STATISTICS* **15**, 651–674 (2006).
 - [36] Tourangeau, R., Rips, L. J. & Rasinski, K. The psychology of survey response (2000).
 - [37] Malhotra, N. Completion time and response order effects in web surveys. *Public opinion quarterly* **72**, 914–934 (2008).
 - [38] Callegaro, M., Yang, Y., Bhola, D. S., Dillman, D. A. & Chin, T.-Y. Response latency as an indicator of optimizing in online questionnaires. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique* **103**, 5–25 (2009).
 - [39] Resnick, P. J. The detection of malingered mental illness. *Behavioral Sciences & the Law* **2**, 21–38 (1984).
 - [40] Rogers, R. & Bender, S. D. *Clinical assessment of malingering and deception* (Guilford Publications, 2020).
 - [41] Dell, N. A., Carbone, J. T., Holzer, K. J. & Vaughn, M. G. Malingering and comorbid psychopathology: Evidence from the 2016-2017 nationwide emergency department sample. *General Hospital Psychiatry* **73**, 121–122 (2021).
 - [42] McDermott, B., Dualan, I. & Scott, C. Malingering in the correctional system: does incentive affect prevalence? *Int'l J L & Psychiatry* **36**, 287–92 (2013).
 - [43] Schmidt, T., Krüger, M. & Ullmann, U. Base rate of probable malingering and its indicators in the assessment of mental disorders-retrospective analysis of a sample of forensic psychological evaluations. *Die Rehabilitation* **59**, 231–236 (2020).
 - [44] DePaulo, B. M., Lindsay, J. J., Malone, B. E., Muhlenbruck, L. & Charlton, K. Cues to deception. *Psychological Bulletin* **129**, 74–118 (2003).
 - [45] Walczyk, J. J., Sewell, N. & DiBenedetto, M. B. A review of approaches to detecting malingering in forensic contexts and promising cognitive load-inducing lie detection techniques. *Frontiers in psychiatry* **9**, 700 (2018).
 - [46] Rogers, R. & Correa, A. A. Determinations of malingering: Evolution from case-based methods to detection strategies. *Psychiatry, Psychology and Law* **15**, 213–223 (2008).
 - [47] Sartori, G., Agosta, S., Zogmaister, C., Ferrara, S. D. & Castiello, U. How to accurately detect autobiographical events. *Psychological science* **19**, 772–780 (2008).
 - [48] Gregg, A. P. When vying reveals lying: The timed antagonistic response alethiometer. *Applied Cognitive Psychology: The Official Journal of the Society for Applied Research in Memory and Cognition* **21**, 621–647 (2007).
 - [49] Monaro, M., Gamberini, L. & Sartori, G. The detection of faked identity using unexpected questions and mouse dynamics. *PloS one* **12**, e0177851 (2017).
 - [50] Sarda-Espinosa, A., Subbiah, S. & Bartz-Beielstein, T. Conditional inference trees for knowledge extraction from motor health condition data. *Engineering Applications of Artificial Intelligence* **62**, 26–37 (2017).
 - [51] Cover, T. M. & Thomas, J. A. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)* (Wiley-Interscience, New York, NY, USA, 2006).
 - [52] Fuglede, B. & Topsoe, F. Jensen-shannon divergence and hilbert space embedding. In *International Symposium on Information Theory, 2004. ISIT 2004. Proceedings.*, 31 (IEEE, 2004).
 - [53] Lundberg, S. M. & Lee, S.-I. A unified approach to interpreting model predictions. In *Advances in neural information processing systems* (Curran Associates, Inc., 2017).
 - [54] Fedotov, A. A., Harremoës, P. & Topsoe, F. Refinements of pinsker's inequality. *IEEE Transactions on Information Theory* **49**, 1491–1498 (2003).

DATA, MATERIALS AND SOFTWARE AVAILABILITY

Software for inferring Q-nets is available as an open-source python package `quasinet`, and can be installed from the standard Python code registry.

ACKNOWLEDGEMENTS

We extend our appreciation to the PS cohort comprising mental health professionals from around Chicago for their uncompensated participation, and the Prolific survey participants, who received nominal compensation for their invaluable contributions to this study.

Funding

This work is funded in part by the Defense Sciences Office of the Defense Advanced Research Projects Agency (Project No. W911NF2010302). The claims made in this study do not necessarily reflect the position or the policy of the sponsors, and no official endorsement should be inferred.

Author Contributions

IC originated the idea, performed analysis, provided funding, developed software and wrote the paper. NC performed analysis, wrote software and and wrote the paper. RG, RL and JE interpreted data and wrote the paper.

Regulatory Approvals

Data collection for the PL and PS cohorts were approved by University of Chicago IRB #IRB24-0310. The third party platform Prolific adheres to rigorous ethical guidelines and privacy policies, ensuring a diverse and reliable participant pool verified through bank-grade ID checks and continuous quality management, aligning with WCAG 2.1 AA standards for accessibility and inclusivity.